

Import libraries/ Dependencies -

```
In [27]: 1 import pandas as pd
2 import numpy as np
3
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.model_selection import train_test_split
9 from sklearn.metrics import precision_score, recall_score, confusion_matrix, precision_recall_curve
10 from sklearn.metrics import accuracy_score, classification_report, f1_score
11 from sklearn.tree import DecisionTreeClassifier
12 from sklearn.neighbors import KNeighborsClassifier
13 from imblearn.combine import SMOTEENN
14
15 import warnings
16 warnings.filterwarnings("ignore")
17
18 from IPython.display import display
19 pd.set_option("display.max_columns", None)
20 pd.set_option("display.max_rows", None)
21 %matplotlib inline
22
23 import json
24 import pickle
```

Reading CSV

```
In [2]: 1 df = pd.read_csv("Tel_churn.csv")
2 df.head()
```

```
Out[2]:
```

tract	PaperlessBilling	MonthlyCharges	TotalCharges	Churn	Tenure1	InternetService_DSL	InternetService_Fiber optic	InternetService
0	1	29.85	29.85	0	1	1	0	
1	0	56.95	1889.50	0	3	1	0	
0	1	53.85	108.15	1	1	1	0	
1	0	42.30	1840.75	0	4	1	0	
0	1	70.70	151.65	1	1	0	1	

```
In [5]: 1 x = df.drop('Churn', axis=1)
2 x.head()
```

```
Out[5]:
```

t	StreamingTV	StreamingMovies	Contract	PaperlessBilling	MonthlyCharges	TotalCharges	Tenure1	InternetService_DSL	Inte
0	0	0	0	1	29.85	29.85	1		1
0	0	0	1	0	56.95	1889.50	3		1
0	0	0	0	1	53.85	108.15	1		1
1	0	0	1	0	42.30	1840.75	4		1
0	0	0	0	1	70.70	151.65	1		0

```
In [7]: 1 y = df["Churn"]
2 y.head()
```

```
Out[7]: 0    0
1    0
2    1
3    0
4    1
Name: Churn, dtype: int64
```

Train Test Split

```
In [35]: 1 # Splitting of dataset for the training and testing
2 x_train, x_test, y_train, y_test = train_test_split(x,y,train_size=0.7,random_state=42)
```

```
In [36]: 1 # Checking for the proper splitting happened or not.
2 print(f"X Train Shape - {x_train.shape}")
3 print(f"X Test Shape - {x_test.shape}")
4 print(f"Y Train Shape - {y_train.shape}")
5 print(f"Y Test Shape - {y_test.shape}")
```

X Train Shape - (4922, 24)

X Test Shape - (2110, 24)

Y Train Shape - (4922,)

Y Test Shape - (2110,)

Model Fitting -

```
In [37]: 1 # Model Training for Logistic regression -
2 lr_model = LogisticRegression()
3 lr_model.fit(x_train,y_train)
```

```
Out[37]: ▾ LogisticRegression
LogisticRegression()
```

```
In [38]: 1 # Model Training for KNN Classifier -
2 knn_model = KNeighborsClassifier()
3 knn_model.fit(x_train,y_train)
```

```
Out[38]: ▾ KNeighborsClassifier
KNeighborsClassifier()
```

```
In [39]: 1 # Model Training for Descision Tree Classifier -
2 dt_model = DecisionTreeClassifier()
3 dt_model.fit(x_train,y_train)
```

```
Out[39]: ▾ DecisionTreeClassifier
DecisionTreeClassifier()
```

Model Evaluation

```
In [40]: 1 # Evaluation of Logistic regression Model based on Training data -
2 y_pred = lr_model.predict(x_train)
3
4 cnf_matrix = confusion_matrix(y_train,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci = precision_score(y_train,y_pred)
8 print(f"Precision - {preci}\n")
9 recal = recall_score(y_train,y_pred)
10 print(f"Recall - {recal}\n")
11 accuracy = accuracy_score(y_train,y_pred)
12 print(f"Accuracy - {accuracy}\n")
13 f1 = f1_score(y_train,y_pred)
14 print(f"F1 Score - {f1}\n")
15
16 cnf_report = classification_report(y_train,y_pred)
17 print(f"Classification report - \n{cnf_report}")
```

Confusion matrix -
[[3225 389]
 [596 712]]

Precision - 0.6466848319709355

Recall - 0.5443425076452599

Accuracy - 0.7998780983340106

F1 Score - 0.5911166459111664

Classification report -

	precision	recall	f1-score	support
0	0.84	0.89	0.87	3614
1	0.65	0.54	0.59	1308
accuracy			0.80	4922
macro avg	0.75	0.72	0.73	4922
weighted avg	0.79	0.80	0.79	4922

```

In [41]: 1 # Evaluation of Logistic regression Model based on Testing data -
2 y_pred = lr_model.predict(x_test)
3
4 cnf_matrix = confusion_matrix(y_test,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci = precision_score(y_test,y_pred)
8 print(f"Precision - {preci}\n")
9 recal = recall_score(y_test,y_pred)
10 print(f"Recall - {recal}\n")
11 accuracy = accuracy_score(y_test,y_pred)
12 print(f"Accuracy - {accuracy}\n")
13 f1 = f1_score(y_test,y_pred)
14 print(f"F1 Score - {f1}\n")
15
16 cnf_report = classification_report(y_test,y_pred)
17 print(f"Classification report - \n{cnf_report}")
18 res = pd.DataFrame([["Logistic Regression",preci,recal,accuracy,f1]],
19                     columns=["Model","Precision","Recall","Accuracy","f1 Score"])

```

Confusion matrix -

```
[[1383 166]
 [ 266 295]]
```

Precision - 0.6399132321041214

Recall - 0.5258467023172906

Accuracy - 0.795260663507109

F1 Score - 0.5772994129158512

Classification report -

	precision	recall	f1-score	support
0	0.84	0.89	0.86	1549
1	0.64	0.53	0.58	561
accuracy			0.80	2110
macro avg	0.74	0.71	0.72	2110
weighted avg	0.79	0.80	0.79	2110

```

In [42]: 1 # Evaluation of KNN classifier Model based on Training data -
2 y_pred = knn_model.predict(x_train)
3
4 cnf_matrix = confusion_matrix(y_train,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci = precision_score(y_train,y_pred)
8 print(f"Precision - {preci}\n")
9 recal = recall_score(y_train,y_pred)
10 print(f"Recall - {recal}\n")
11 accuracy = accuracy_score(y_train,y_pred)
12 print(f"Accuracy - {accuracy}\n")
13 f1 = f1_score(y_train,y_pred)
14 print(f"F1 Score - {f1}\n")
15
16 cnf_report = classification_report(y_train,y_pred)
17 print(f"Classification report - \n{cnf_report}")

```

Confusion matrix -

```
[[3323 291]
 [ 543 765]]
```

Precision - 0.7244318181818182

Recall - 0.5848623853211009

Accuracy - 0.8305566842746851

F1 Score - 0.6472081218274112

Classification report -

	precision	recall	f1-score	support
0	0.86	0.92	0.89	3614
1	0.72	0.58	0.65	1308
accuracy			0.83	4922
macro avg	0.79	0.75	0.77	4922
weighted avg	0.82	0.83	0.82	4922

```

In [43]: 1 # Evaluation of KNN classifier Model based on Testing data -
2 y_pred = knn_model.predict(x_test)
3
4 cnf_matrix = confusion_matrix(y_test,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci1 = precision_score(y_test,y_pred)
8 print(f"Precision - {preci1}\n")
9 recal1 = recall_score(y_test,y_pred)
10 print(f"Recall - {recal1}\n")
11 accuracy1 = accuracy_score(y_test,y_pred)
12 print(f"Accuracy - {accuracy1}\n")
13 f11 = f1_score(y_test,y_pred)
14 print(f"F1 Score - {f11}\n")
15
16 cnf_report = classification_report(y_test,y_pred)
17 print(f"Classification report - \n{cnf_report}")
18 res1 = pd.DataFrame([[ "KNN Classifier",preci1,recal1,accuracy1,f11]],
19                      columns=["Model", "Precision", "Recall", "Accuracy", "f1 Score"])

```

Confusion matrix -

```
[[1379  170]
 [ 312  249]]
```

Precision - 0.594272076372315

Recall - 0.44385026737967914

Accuracy - 0.771563981042654

F1 Score - 0.5081632653061224

Classification report -

	precision	recall	f1-score	support
0	0.82	0.89	0.85	1549
1	0.59	0.44	0.51	561
accuracy			0.77	2110
macro avg	0.70	0.67	0.68	2110
weighted avg	0.76	0.77	0.76	2110

```
In [44]: 1 # Evaluation of Descision Tree classifier Model based on Training data -
2 y_pred = dt_model.predict(x_train)
3
4 cnf_matrix = confusion_matrix(y_train,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci = precision_score(y_train,y_pred)
8 print(f"Precision - {preci}\n")
9 recal = recall_score(y_train,y_pred)
10 print(f"Recall - {recal}\n")
11 accuracy = accuracy_score(y_train,y_pred)
12 print(f"Accuracy - {accuracy}\n")
13 f1 = f1_score(y_train,y_pred)
14 print(f"F1 Score - {f1}\n")
15
16 cnf_report = classification_report(y_train,y_pred)
17 print(f"Classification report - \n{cnf_report}")
```

Confusion matrix -

```
[[3613   1]
 [   6 1302]]
```

Precision - 0.9992325402916347

Recall - 0.9954128440366973

Accuracy - 0.9985778138967899

F1 Score - 0.9973190348525469

Classification report -

	precision	recall	f1-score	support
0	1.00	1.00	1.00	3614
1	1.00	1.00	1.00	1308
accuracy			1.00	4922
macro avg	1.00	1.00	1.00	4922
weighted avg	1.00	1.00	1.00	4922

In [45]:

```

1 # Evaluation of Descision Tree classifier Model based on Testing data -
2 y_pred = dt_model.predict(x_test)
3
4 cnf_matrix = confusion_matrix(y_test,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci2 = precision_score(y_test,y_pred)
8 print(f"Precision - {preci2}\n")
9 recal2 = recall_score(y_test,y_pred)
10 print(f"Recall - {recal2}\n")
11 accuracy2 = accuracy_score(y_test,y_pred)
12 print(f"Accuracy - {accuracy2}\n")
13 f12 = f1_score(y_test,y_pred)
14 print(f"F1 Score - {f12}\n")
15
16 cnf_report = classification_report(y_test,y_pred)
17 print(f"Classification report - \n{cnf_report}")
18 res2 = pd.DataFrame([["Descision Tree Classifier",preci2,recal2,accuracy2,f12]],
19                     columns=["Model","Precision","Recall","Accuracy","f1 Score"])

```

Confusion matrix -

```
[[1234  315]
 [ 260  301]]
```

Precision - 0.48863636363636365

Recall - 0.5365418894830659

Accuracy - 0.7274881516587678

F1 Score - 0.5114698385726423

Classification report -

	precision	recall	f1-score	support
0	0.83	0.80	0.81	1549
1	0.49	0.54	0.51	561
accuracy			0.73	2110
macro avg	0.66	0.67	0.66	2110
weighted avg	0.74	0.73	0.73	2110

In [46]:

```

1 # Creating new dataframe by joining 3 dataframes of evaluation values of different models
2 result = pd.concat([res,res1,res2],ignore_index=True)
3 result

```

Out[46]:

	Model	Precision	Recall	Accuracy	f1 Score
0	Logistic Regression	0.639913	0.525847	0.795261	0.577299
1	KNN Classifier	0.594272	0.443850	0.771564	0.508163
2	Descision Tree Classifier	0.488636	0.536542	0.727488	0.511470

- As you can see that the accuracy is quite low, and as it's an imbalanced dataset, we shouldn't consider Accuracy as our metrics to measure the model, as Accuracy is cursed in imbalanced datasets.
- Hence, we need to check recall, precision & f1 score for the minority class, and it's quite evident that the precision, recall & f1 score is too low for Class 1, i.e. churned customers.
- Hence, moving ahead to call SMOTEENN (UpSampling + ENN)

Resampling SMOTEENN

In [47]:

```

1 sm = SMOTEENN()
2 X_resampled, y_resampled = sm.fit_resample(x,y)

```

In [48]:

```

1 # Splitting of dataset for the training and testing
2 xr_train,xr_test,yr_train,yr_test=train_test_split(X_resampled, y_resampled,test_size=0.2)

```



```
In [49]: 1 # Checking for the proper splitting happened or not.  
2 print(f"X Train Shape - {xr_train.shape}")  
3 print(f"X Test Shape - {xr_test.shape}")  
4 print(f"Y Train Shape - {yr_train.shape}")  
5 print(f"Y Test Shape - {yr_test.shape}")
```

X Train Shape - (4615, 24)
X Test Shape - (1154, 24)
Y Train Shape - (4615,)
Y Test Shape - (1154,)

Model Training After SMOTE

```
In [50]: 1 # Model Training for logistic regression after sampling -  
2 lr_model_sm = LogisticRegression()  
3 lr_model_sm.fit(xr_train,yr_train)
```

```
Out[50]: ▾ LogisticRegression  
LogisticRegression()
```

```
In [51]: 1 # Model Training for KNN Classifier after sampling-  
2 knn_model_sm = KNeighborsClassifier()  
3 knn_model_sm.fit(xr_train,yr_train)
```

```
Out[51]: ▾ KNeighborsClassifier  
KNeighborsClassifier()
```

```
In [52]: 1 # Model Training for Descision Tree Classifier after sampling -  
2 dt_model_sm = DecisionTreeClassifier()  
3 dt_model_sm.fit(xr_train,yr_train)
```

```
Out[52]: ▾ DecisionTreeClassifier  
DecisionTreeClassifier()
```

Model Evaluation After SMOTE

In [60]:

```

1 # Evaluation of Logistic regression Model based on Training data after sampling-
2 y_pred = lr_model_sm.predict(xr_train)
3
4 cnf_matrix = confusion_matrix(yr_train,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci = precision_score(yr_train,y_pred)
8 print(f"Precision - {preci}\n")
9 recal = recall_score(yr_train,y_pred)
10 print(f"Recall - {recal}\n")
11 accuracy = accuracy_score(yr_train,y_pred)
12 print(f"Accuracy - {accuracy}\n")
13 f1 = f1_score(yr_train,y_pred)
14 print(f"F1 Score - {f1}\n")
15
16 cnf_report = classification_report(yr_train,y_pred)
17 print(f"Classification report - \n{cnf_report}")

```

Confusion matrix -

```
[[1948 167]
 [ 151 2349]]
```

Precision - 0.9336248012718601

Recall - 0.9396

Accuracy - 0.9310942578548213

F1 Score - 0.9366028708133971

Classification report -

	precision	recall	f1-score	support
0	0.93	0.92	0.92	2115
1	0.93	0.94	0.94	2500
accuracy			0.93	4615
macro avg	0.93	0.93	0.93	4615
weighted avg	0.93	0.93	0.93	4615

```

In [61]: 1 # Evaluation of Logistic regression Model based on Testing data after sampling-
2 y_pred = lr_model_sm.predict(xr_test)
3
4 cnf_matrix = confusion_matrix(yr_test,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci = precision_score(yr_test,y_pred)
8 print(f"Precision - {preci}\n")
9 recal = recall_score(yr_test,y_pred)
10 print(f"Recall - {recal}\n")
11 accuracy = accuracy_score(yr_test,y_pred)
12 print(f"Accuracy - {accuracy}\n")
13 f1 = f1_score(yr_test,y_pred)
14 print(f"F1 Score - {f1}\n")
15
16 cnf_report = classification_report(yr_test,y_pred)
17 print(f"Classification report - \n{cnf_report}")
18 res11 = pd.DataFrame([["Logistic Regression",preci,recal,accuracy,f1]],
19                       columns=["Model","Precision","Recall","Accuracy","f1 Score"])

```

Confusion matrix -

```
[[472  55]
 [ 39 588]]
```

Precision - 0.9144634525660964

Recall - 0.937799043062201

Accuracy - 0.9185441941074524

F1 Score - 0.925984251968504

Classification report -

	precision	recall	f1-score	support
0	0.92	0.90	0.91	527
1	0.91	0.94	0.93	627
accuracy			0.92	1154
macro avg	0.92	0.92	0.92	1154
weighted avg	0.92	0.92	0.92	1154

```
In [62]: 1 # Evaluation of KNN classifier Model based on Training data after sampling-
2 y_pred = knn_model_sm.predict(xr_train)
3
4 cnf_matrix = confusion_matrix(yr_train,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci = precision_score(yr_train,y_pred)
8 print(f"Precision - {preci}\n")
9 recal = recall_score(yr_train,y_pred)
10 print(f"Recall - {recal}\n")
11 accuracy = accuracy_score(yr_train,y_pred)
12 print(f"Accuracy - {accuracy}\n")
13 f1 = f1_score(yr_train,y_pred)
14 print(f"F1 Score - {f1}\n")
15
16 cnf_report = classification_report(yr_train,y_pred)
17 print(f"Classification report - \n{cnf_report}")
```

Confusion matrix -

```
[[2024  91]
 [ 31 2469]]
```

Precision - 0.964453125

Recall - 0.9876

Accuracy - 0.9735644637053088

F1 Score - 0.9758893280632412

Classification report -

	precision	recall	f1-score	support
0	0.98	0.96	0.97	2115
1	0.96	0.99	0.98	2500
accuracy			0.97	4615
macro avg	0.97	0.97	0.97	4615
weighted avg	0.97	0.97	0.97	4615

```

In [64]: 1 # Evaluation of KNN classifier Model based on Testing data after sampling-
2 y_pred = knn_model_sm.predict(xr_test)
3
4 cnf_matrix = confusion_matrix(yr_test,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci1 = precision_score(yr_test,y_pred)
8 print(f"Precision - {preci1}\n")
9 recal1 = recall_score(yr_test,y_pred)
10 print(f"Recall - {recal1}\n")
11 accuracy1 = accuracy_score(yr_test,y_pred)
12 print(f"Accuracy - {accuracy1}\n")
13 f11 = f1_score(yr_test,y_pred)
14 print(f"F1 Score - {f11}\n")
15
16 cnf_report = classification_report(yr_test,y_pred)
17 print(f"Classification report - \n{cnf_report}")
18 res22 = pd.DataFrame([["KNN Classifier",preci1,recal1,accuracy1,f11]],
19                      columns=["Model","Precision","Recall","Accuracy","f1 Score"])

```

Confusion matrix -

```
[[485  42]
 [ 13 614]]
```

Precision - 0.9359756097560976

Recall - 0.9792663476874003

Accuracy - 0.9523396880415944

F1 Score - 0.9571317225253314

Classification report -

	precision	recall	f1-score	support
0	0.97	0.92	0.95	527
1	0.94	0.98	0.96	627
accuracy			0.95	1154
macro avg	0.95	0.95	0.95	1154
weighted avg	0.95	0.95	0.95	1154

```
In [65]: 1 # Evaluation of Descision Tree classifier Model based on Training data after sampling-
2 y_pred = dt_model_sm.predict(xr_train)
3
4 cnf_matrix = confusion_matrix(yr_train,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci = precision_score(yr_train,y_pred)
8 print(f"Precision - {preci}\n")
9 recal = recall_score(yr_train,y_pred)
10 print(f"Recall - {recal}\n")
11 accuracy = accuracy_score(yr_train,y_pred)
12 print(f"Accuracy - {accuracy}\n")
13 f1 = f1_score(yr_train,y_pred)
14 print(f"F1 Score - {f1}\n")
15
16 cnf_report = classification_report(yr_train,y_pred)
17 print(f"Classification report - \n{cnf_report}")
```

Confusion matrix -
[[2115 0]
[0 2500]]

Precision - 1.0

Recall - 1.0

Accuracy - 1.0

F1 Score - 1.0

Classification report -

	precision	recall	f1-score	support
0	1.00	1.00	1.00	2115
1	1.00	1.00	1.00	2500
accuracy			1.00	4615
macro avg	1.00	1.00	1.00	4615
weighted avg	1.00	1.00	1.00	4615

```
In [67]: 1 # Evaluation of Descision Tree classifier Model based on Testing data after sampling -
2 y_pred = dt_model_sm.predict(xr_test)
3
4 cnf_matrix = confusion_matrix(yr_test,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci2 = precision_score(yr_test,y_pred)
8 print(f"Precision - {preci2}\n")
9 recal2 = recall_score(yr_test,y_pred)
10 print(f"Recall - {recal2}\n")
11 accuracy2 = accuracy_score(yr_test,y_pred)
12 print(f"Accuracy - {accuracy2}\n")
13 f12 = f1_score(yr_test,y_pred)
14 print(f"F1 Score - {f12}\n")
15
16 cnf_report = classification_report(yr_test,y_pred)
17 print(f"Classification report - \n{cnf_report}")
18 res33 = pd.DataFrame([["Descision Tree Classifier",preci2,recal2,accuracy2,f12]],
19                      columns=["Model","Precision","Recall","Accuracy","f1 Score"])
```

Confusion matrix -

```
[[479  48]
 [ 33 594]]
```

Precision - 0.9252336448598131

Recall - 0.9473684210526315

Accuracy - 0.9298093587521664

F1 Score - 0.9361702127659575

Classification report -

	precision	recall	f1-score	support
0	0.94	0.91	0.92	527
1	0.93	0.95	0.94	627
accuracy			0.93	1154
macro avg	0.93	0.93	0.93	1154
weighted avg	0.93	0.93	0.93	1154

```
In [68]: 1 # Creating new dataframe by joining 3 dataframes of evaluation values of different models after s
2 result1 = pd.concat([res11,res22,res33],ignore_index=True)
3 result1
```

Out[68]:

	Model	Precision	Recall	Accuracy	f1 Score
0	Logistic Regression	0.914463	0.937799	0.918544	0.925984
1	KNN Classifier	0.935976	0.979266	0.952340	0.957132
2	Descision Tree Classifier	0.925234	0.947368	0.929809	0.936170

As observed, for the context of our project, we are aiming for f1Score. KNN Classifier presents the highest f1Score score. Lets assume we are aiming more for recall and precision, we could then, take advantage of f1-score, as its the harmonic average between precision and recall. Thus, the algorithm that satisfies this need is KNN Classifier . but KNN Classifier is lazy learner so we use either Logistic Regression or Descision tree .

Now, we can see quite better results and a very good recall, precision & f1 score for minority class. We can try pruning on Descision Tree because there is scope of improvements in DT due to overfitting.

Pruning - for reducing Overfitting

```
In [69]: 1 # Model training of Dt for ccp alpha values after sampling -
2 dt_clf_ccp = DecisionTreeClassifier(random_state=10,ccp_alpha=0.0)
3 dt_clf_ccp.fit(xr_train,yr_train)
```

```
Out[69]: DecisionTreeClassifier
DecisionTreeClassifier(random_state=10)
```

```
In [70]: 1 result1 = dt_clf_ccp.cost_complexity_pruning_path(xr_train,yr_train)
2 ccp_list = result1["ccp_alphas"]
3 ccp_list
```

```
Out[70]: array([0.00000000e+00, 1.08239277e-04, 1.35427952e-04, 1.39106242e-04,
1.39297322e-04, 1.39297322e-04, 1.41922158e-04, 1.43040242e-04,
1.44456482e-04, 1.44597034e-04, 1.89599133e-04, 1.89599133e-04,
1.90082805e-04, 1.92608643e-04, 1.95016251e-04, 1.96943732e-04,
2.00016668e-04, 2.01207243e-04, 2.01207243e-04, 2.03141928e-04,
2.05280265e-04, 2.05850488e-04, 2.06366404e-04, 2.06366404e-04,
2.07464097e-04, 2.11974186e-04, 2.11974186e-04, 2.14580989e-04,
2.17325338e-04, 2.17505499e-04, 2.18146834e-04, 2.50391236e-04,
2.56120377e-04, 2.70855905e-04, 2.79339037e-04, 2.82906026e-04,
2.84215193e-04, 2.85572931e-04, 2.88912965e-04, 2.88912965e-04,
2.88912965e-04, 2.88912965e-04, 3.03358613e-04, 3.09549605e-04,
3.15008330e-04, 3.25027086e-04, 3.25027086e-04, 3.25027086e-04,
3.46695558e-04, 3.46695558e-04, 3.46695558e-04, 3.46695558e-04,
3.48759222e-04, 3.61141206e-04, 3.61141206e-04, 3.61141206e-04,
3.64150716e-04, 3.71459526e-04, 3.71601848e-04, 3.75211643e-04,
3.79198267e-04, 3.85217287e-04, 3.87510247e-04, 3.90032503e-04,
3.90032503e-04, 3.93972225e-04, 4.11276103e-04, 4.16061297e-04,
4.36944598e-04, 4.44481485e-04, 4.60514660e-04, 4.83026363e-04,
4.85373781e-04, 4.87540628e-04, 4.95161197e-04, 5.08845690e-04,
5.20043337e-04, 5.20043337e-04, 5.39476173e-04, 5.41711809e-04,
5.90958337e-04, 6.29969313e-04, 6.42533937e-04, 6.48787009e-04,
6.50054171e-04, 6.50054171e-04, 6.50054171e-04, 6.58900307e-04,
6.82249586e-04, 7.07836764e-04, 7.45806583e-04, 7.47170698e-04,
7.97873831e-04, 8.12156900e-04, 8.13771518e-04, 8.53196100e-04,
9.28555608e-04, 9.61080090e-04, 9.93998177e-04, 1.00408825e-03,
1.02467798e-03, 1.08924282e-03, 1.10368342e-03, 1.15148427e-03,
1.17039278e-03, 1.20380402e-03, 1.20432673e-03, 1.25065764e-03,
1.27595078e-03, 1.31480413e-03, 1.55072202e-03, 1.66529123e-03,
2.40914897e-03, 3.00027287e-03, 3.91787900e-03, 4.26936819e-03,
4.27113288e-03, 4.42286692e-03, 4.95869592e-03, 5.46892719e-03,
8.64953652e-03, 1.39933919e-02, 5.00585556e-02, 2.63838098e-01])
```

```
In [71]: 1 test_acc_list = []
2 train_acc_list = []
3 for ccp in ccp_list:
4     dt_clf_ccp = DecisionTreeClassifier(random_state=10,ccp_alpha=ccp)
5     dt_clf_ccp.fit(xr_train,yr_train)
6     train_acc_list.append(dt_clf_ccp.score(xr_train,yr_train))
7     test_acc_list.append(dt_clf_ccp.score(xr_test,yr_test))
```

```
In [72]: 1 index = np.argmax(test_acc_list)
2 index
```

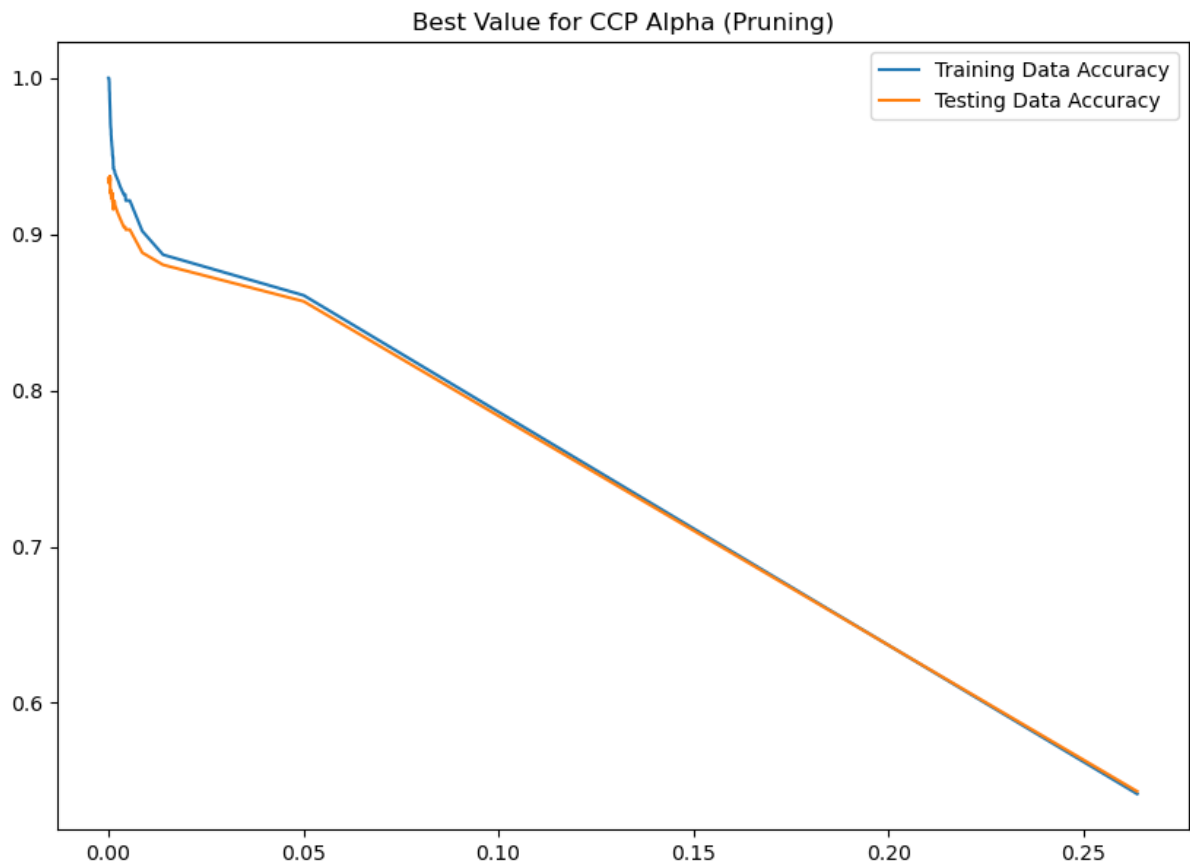
```
Out[72]: 31
```

```
In [73]: 1 # best value of ccp alpha
2 cc_val = ccp_list[index]
3 print("Best Value of CCP Alpha - ",cc_val)
```

```
Best Value of CCP Alpha - 0.0002503912363067293
```



```
In [74]: 1 # Plot of Accuracies to take best value of ccp_alpha from graph
2 fig,ax = plt.subplots()
3 ax.figure.set_size_inches(10,7)
4 ax.plot(ccp_list,train_acc_list,label="Training Data Accuracy")
5 ax.plot(ccp_list,test_acc_list,label="Testing Data Accuracy")
6 ax.legend()
7 plt.title("Best Value for CCP Alpha (Pruning)")
8 plt.show()
```



```
In [75]: 1 dt_clf_ccp_sm = DecisionTreeClassifier(random_state=10,ccp_alpha=cc_val)
2 dt_clf_ccp_sm.fit(xr_train,yr_train)
```

```
Out[75]: DecisionTreeClassifier
DecisionTreeClassifier(ccp_alpha=0.0002503912363067293, random_state=10)
```

```
In [76]: 1 # Evaluation of Descision Tree classifier Model based on Training data after sampling and Pruning
2 y_pred = dt_clf_ccp_sm.predict(xr_train)
3
4 cnf_matrix = confusion_matrix(yr_train,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci = precision_score(yr_train,y_pred)
8 print(f"Precision - {preci}\n")
9 recal = recall_score(yr_train,y_pred)
10 print(f"Recall - {recal}\n")
11 accuracy = accuracy_score(yr_train,y_pred)
12 print(f"Accuracy - {accuracy}\n")
13 f1 = f1_score(yr_train,y_pred)
14 print(f"F1 Score - {f1}\n")
15
16 cnf_report = classification_report(yr_train,y_pred)
17 print(f"Classification report - \n{cnf_report}")
```

Confusion matrix -

```
[[2092  23]
 [ 17 2483]]
```

Precision - 0.9908220271348763

Recall - 0.9932

Accuracy - 0.991332611050921

F1 Score - 0.9920095884938074

Classification report -

	precision	recall	f1-score	support
0	0.99	0.99	0.99	2115
1	0.99	0.99	0.99	2500
accuracy			0.99	4615
macro avg	0.99	0.99	0.99	4615
weighted avg	0.99	0.99	0.99	4615

```
In [77]: 1 # Evaluation of Descision Tree classifier Model based on Testing data after sampling and Pruning -
2 y_pred = dt_clf_ccp_sm.predict(xr_test)
3
4 cnf_matrix = confusion_matrix(yr_test,y_pred)
5 print(f"Confusion matrix - \n{cnf_matrix}\n")
6
7 preci2 = precision_score(yr_test,y_pred)
8 print(f"Precision - {preci2}\n")
9 recal2 = recall_score(yr_test,y_pred)
10 print(f"Recall - {recal2}\n")
11 accuracy2 = accuracy_score(yr_test,y_pred)
12 print(f"Accuracy - {accuracy2}\n")
13 f12 = f1_score(yr_test,y_pred)
14 print(f"F1 Score - {f12}\n")
15
16 cnf_report = classification_report(yr_test,y_pred)
17 print(f"Classification report - \n{cnf_report}")
18 res33 = pd.DataFrame([["Descision Tree Classifier",preci2,recal2,accuracy2,f12]],
19                      columns=["Model","Precision","Recall","Accuracy","f1 Score"])
```

Confusion matrix -

```
[[481  46]
 [ 27 600]]
```

Precision - 0.9287925696594427

Recall - 0.9569377990430622

Accuracy - 0.9367417677642981

F1 Score - 0.9426551453260016

Classification report -

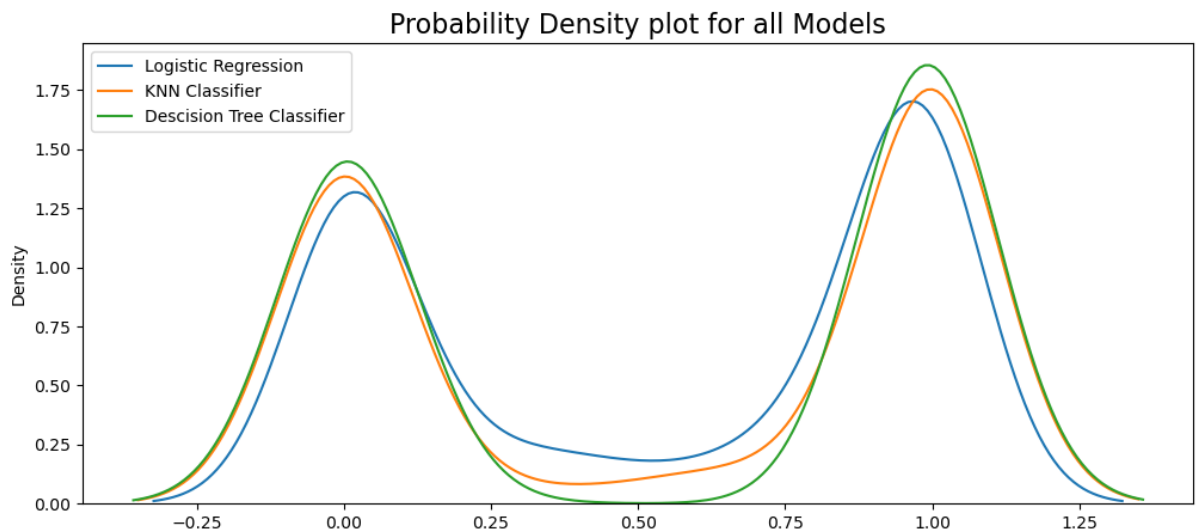
	precision	recall	f1-score	support
0	0.95	0.91	0.93	527
1	0.93	0.96	0.94	627
accuracy			0.94	1154
macro avg	0.94	0.93	0.94	1154
weighted avg	0.94	0.94	0.94	1154

There is quite improvement of accuracy and f1score in the Descision Tree model now we compare predicting probability of classes prediction and then we can finalize the Best model.

Probability Distribution

```
In [79]: 1 # Storing predicted probabilities for class 1
2 y_pred_lr_prob = lr_model_sm.predict_proba(xr_test)[: ,1]
3 y_pred_knn_prob = knn_model_sm.predict_proba(xr_test)[: ,1]
4 y_pred_dt_prob = dt_clf_ccp_sm.predict_proba(xr_test)[: ,1]
```

```
In [80]: 1 # Plotting kdeplot of Probability density for all models
2 plt.figure(figsize=(12,5))
3 sns.kdeplot(y_pred_lr_prob,label="Logistic Regression")
4 sns.kdeplot(y_pred_knn_prob,label="KNN Classifier")
5 sns.kdeplot(y_pred_dt_prob,label="Decision Tree Classifier")
6 plt.title("Probability Density plot for all Models",fontsize=16)
7 plt.legend()
8 plt.show()
```



As observed, in general all the algorithms presents most probabilities concentrated around 1 . However, DescisionTreeClassifier presents the largest amount of probabilities concentrated around 1 and 0 , while DescisionTreeClassifier presenting the least with a moderate distributed probabilities around other values. The term DescisionTreeClassifier suits best for the context and challenge of our project.

Here, We have got the best values of accuracy, precision, recall and f1 score. Descision Tree model after pruning and avoiding overfitting. Descision Tree Model Accuracies -

- Training -
 - Precision - 0.9908220271348763
 - Recall - 0.9932
 - Accuracy - 0.991332611050921
 - F1 Score - 0.9920095884938074
- Testing -
 - Precision - 0.9287925696594427
 - Recall - 0.9569377990430622
 - Accuracy - 0.9367417677642981
 - F1 Score - 0.9426551453260016

User Input Testing

```
In [82]: 1 x.columns
```

```
Out[82]: Index(['gender', 'SeniorCitizen', 'Partner', 'Dependents', 'PhoneService',
               'MultipleLines', 'OnlineSecurity', 'OnlineBackup', 'DeviceProtection',
               'TechSupport', 'StreamingTV', 'StreamingMovies', 'Contract',
               'PaperlessBilling', 'MonthlyCharges', 'TotalCharges', 'Tenure1',
               'InternetService_DSL', 'InternetService_Fiber optic',
               'InternetService_No', 'PaymentMethod_Bank transfer (automatic)',
               'PaymentMethod_Credit card (automatic)',
               'PaymentMethod_Electronic check', 'PaymentMethod_Mailed check'],
              dtype='object')
```

```

In [92]: 1 # Writting code to create dictionary of all encoded Features -
2 gender_val = {'Male': 1, 'Female': 0}
3 phone_service_val = {'Yes': 1, 'No': 0}
4 dependents_val = {'No': 0, 'Yes': 1}
5 partner_val = {'No': 0, 'Yes': 1}
6 multiple_lines_val = {'No': 0, 'Yes': 1, 'No phone service': 2}
7 online_security_val = {'No': 0, 'Yes': 1, 'No internet service': 2}
8 online_backup_val = {'No': 0, 'Yes': 1, 'No internet service': 2}
9 device_protection_val = {'No': 0, 'Yes': 1, 'No internet service': 2}
10 tech_support_val = {'No': 0, 'Yes': 1, 'No internet service': 2}
11 streaming_tv_val = {'No': 0, 'Yes': 1, 'No internet service': 2}
12 streaming_movies_val = {'No': 0, 'Yes': 1, 'No internet service': 2}
13 contract_val = {'Month-to-month': 0, 'Two year': 2, 'One year': 1}
14 paper_less_billing_val = {'Yes': 1, 'No': 0}
15 encoded = {"gender_val" : gender_val,
16            "phone_service_val" : phone_service_val,
17            "dependents_val" : dependents_val,
18            "partner_val" : partner_val,
19            "multiple_lines_val" : multiple_lines_val,
20            "online_security_val" : online_security_val,
21            "online_backup_val" : online_backup_val,
22            "device_protection_val" : device_protection_val,
23            "tech_support_val" : tech_support_val,
24            "streaming_tv_val" : streaming_tv_val,
25            "streaming_movies_val" : streaming_movies_val,
26            "contract_val" : contract_val,
27            "paper_less_billing_val" : paper_less_billing_val,
28            "columns":list(x.columns)}
29 encoded

```

```

Out[92]: {'gender_val': {'Male': 1, 'Female': 0},
'phone_service_val': {'Yes': 1, 'No': 0},
'dependents_val': {'No': 0, 'Yes': 1},
'partner_val': {'No': 0, 'Yes': 1},
'multiple_lines_val': {'No': 0, 'Yes': 1, 'No phone service': 2},
'online_security_val': {'No': 0, 'Yes': 1, 'No internet service': 2},
'online_backup_val': {'No': 0, 'Yes': 1, 'No internet service': 2},
'device_protection_val': {'No': 0, 'Yes': 1, 'No internet service': 2},
'tech_support_val': {'No': 0, 'Yes': 1, 'No internet service': 2},
'streaming_tv_val': {'No': 0, 'Yes': 1, 'No internet service': 2},
'streaming_movies_val': {'No': 0, 'Yes': 1, 'No internet service': 2},
'contract_val': {'Month-to-month': 0, 'Two year': 2, 'One year': 1},
'paper_less_billing_val': {'Yes': 1, 'No': 0},
'columns': ['gender',
'SeniorCitizen',
'Partner',
'Dependents',
'PhoneService',
'MultipleLines',
'OnlineSecurity',
'OnlineBackup',
'DeviceProtection',
'TechSupport',
'StreamingTV',
'StreamingMovies',
'Contract',
'PaperlessBilling',
'MonthlyCharges',
'TotalCharges',
'Tenure1',
'InternetService_DSL',
'InternetService_Fiber optic',
'InternetService_No',
'PaymentMethod_Bank transfer (automatic)',
'PaymentMethod_Credit card (automatic)',
'PaymentMethod_Electronic check',
'PaymentMethod_Mailed check']}

```

```

In [93]: 1 # Creating json file of encoded so we can further use it.
2 with open("encoded.json","w") as f:
3     json.dump(encoded,f)

```

In [105]:

```
1 # Opening original Df for the taking random entry for user input testing -
2 df_ut = pd.read_csv("Telco-Customer-Churn.csv").drop(["customerID"],axis=1)
3 df_ut.head()
```

Out[105]:

	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	OnlineSecurity	OnlineBa
0	Female	0	Yes	No	1	No	No phone service	DSL	No	
1	Male	0	No	No	34	Yes	No	DSL	Yes	
2	Male	0	No	No	2	Yes	No	DSL	Yes	
3	Male	0	No	No	45	No	No phone service	DSL	Yes	
4	Female	0	No	No	2	Yes	No	Fiber optic	No	

In [109]:

```
1 # Checking for the 1st entry in the data.
2 df_ut.iloc[2]
```

Out[109]:

genderMale
SeniorCitizen0
PartnerNo
DependentsNo
tenure2
PhoneServiceYes
MultipleLinesNo
InternetServiceDSL
OnlineSecurityYes
OnlineBackupYes
DeviceProtectionNo
TechSupportNo
StreamingTVNo
StreamingMoviesNo
ContractMonth-to-month
PaperlessBillingYes
PaymentMethodMailed check
MonthlyCharges53.85
TotalCharges108.15
ChurnYes
Name: 2, dtype: object

In [114]:

```
1 # Creating variables and assigning known values for user input testing -
2 gender = "Male"
3 SeniorCitizen = 0
4 Partner = "No"
5 Dependents = "No"
6 tenure = 2
7 PhoneService = "Yes"
8 MultipleLines = "No"
9 InternetService = "DSL" #One-Hot-Encoded
10 OnlineSecurity = "Yes"
11 OnlineBackup = "Yes"
12 DeviceProtection = "No"
13 TechSupport = "No"
14 StreamingTV = "No"
15 StreamingMovies = "No"
16 Contract = "Month-to-month"
17 PaperlessBilling = "Yes"
18 PaymentMethod = "Mailed check" #One-Hot-Encoded
19 MonthlyCharges = 53.85
20 TotalCharges = 108.15
```

In [115]: 1 x.columns

Out[115]: Index(['gender', 'SeniorCitizen', 'Partner', 'Dependents', 'PhoneService', 'MultipleLines', 'OnlineSecurity', 'OnlineBackup', 'DeviceProtection', 'TechSupport', 'StreamingTV', 'StreamingMovies', 'Contract', 'PaperlessBilling', 'MonthlyCharges', 'TotalCharges', 'Tenure1', 'InternetService_DSL', 'InternetService_Fiber optic', 'InternetService_No', 'PaymentMethod_Bank transfer (automatic)', 'PaymentMethod_Credit card (automatic)', 'PaymentMethod_Electronic check', 'PaymentMethod_Mailed check'], dtype='object')

```
In [116]: 1 # Creating array of values so we can pass the values to the our model -
2 column_names = x.columns
3
4 gender_1 = gender_val[gender]
5 PhoneService_1 = phone_service_val[PhoneService]
6 Dependents_1 = dependents_val[Dependents]
7 Partner_1 = partner_val[Partner]
8 MultipleLines_1 = multiple_lines_val[MultipleLines]
9 OnlineSecurity_1 = online_security_val[OnlineSecurity]
10 OnlineBackup_1 = online_backup_val[OnlineBackup]
11 DeviceProtection_1 = device_protection_val[DeviceProtection]
12 TechSupport_1 = tech_support_val[TechSupport]
13 StreamingTV_1 = streaming_tv_val[StreamingTV]
14 StreamingMovies_1 = streaming_movies_val[StreamingMovies]
15 Contract_1 = contract_val[Contract]
16 PaperlessBilling_1 = paper_less_billing_val[PaperlessBilling]
17
18 # If-Elif-Else block to encode tenure Feature.
19 if tenure > 0 and tenure <= 12:
20     Tenure1 = 1
21 elif tenure > 12 and tenure <= 24:
22     Tenure1 = 2
23 elif tenure > 24 and tenure <= 36:
24     Tenure1 = 3
25 elif tenure > 36 and tenure <= 48:
26     Tenure1 = 4
27 elif tenure > 48 and tenure <= 60:
28     Tenure1 = 5
29 else:
30     Tenure1 = 6
31
32 array = np.zeros(len(x.columns),dtype=float)
33 array[0] = gender_1
34 array[1] = SeniorCitizen
35 array[2] = Partner_1
36 array[3] = Dependents_1
37 array[4] = PhoneService_1
38 array[5] = MultipleLines_1
39 array[6] = OnlineSecurity_1
40 array[7] = OnlineBackup_1
41 array[8] = DeviceProtection_1
42 array[9] = TechSupport_1
43 array[10] = StreamingTV_1
44 array[11] = StreamingMovies_1
45 array[12] = Contract_1
46 array[13] = PaperlessBilling_1
47 array[14] = MonthlyCharges
48 array[15] = TotalCharges
49 array[16] = Tenure1
50
51 InternetService_1 = "InternetService_" + InternetService
52 InternetService_index = np.where(column_names == InternetService_1)[0][0]
53 array[InternetService_index] = 1
54
55 PaymentMethod_1 = "PaymentMethod_" + PaymentMethod
56 PaymentMethod_index = np.where(column_names == PaymentMethod_1)[0][0]
57 array[PaymentMethod_index] = 1
```

```
In [118]: 1 # predicting values based on user input from the best model.
2 print(array)
3 Predict_Customer_Churn = np.around(dt_clf_ccp_sm.predict([array])[0])
4 print("\nPredicted Customer Churn - ",Predict_Customer_Churn)
5 if Predict_Customer_Churn == 0:
6     print("\nPredicted Customer will not Churn.")
7 else:
8     print("\nPredicted Customer will Churn.")
```

```
[ 1.    0.    0.    0.    1.    0.    1.    1.    0.    0.
  0.    0.    0.    1.   53.85 108.15    1.    1.    0.    0.
  0.    0.    0.    1. ]
```

Predicted Customer Churn - 1

Predicted Customer will Churn.

```
In [119]: 1 # Creating pickle file for Descision Tree model after pruning and sampling -
2 with open("dt_model.pkl","wb") as f1:
3     pickle.dump(dt_clf_ccp_sm,f1)
```